

Construction of a variability map for eukaryotic large subunit ribosomal RNA

Abdelghani Ben Ali, Jan Wuyts, Rupert De Wachter, Axel Meyer¹ and Yves Van de Peer*

Department of Biochemistry, University of Antwerp (UIA), Universiteitsplein 1, B-2610 Antwerpen, Belgium and

¹Department of Biology, University of Konstanz, D-78457 Konstanz, Germany

Received May 14, 1999; Accepted May 27, 1999

ABSTRACT

In this paper, we present a variability map of the eukaryotic large subunit ribosomal RNA, showing the distribution of variable and conserved sites in this molecule. The variability of each site in this map is indicated by means of a colored dot. Construction of the variability map was based on the substitution rate calibration (SRC) method, in which the substitution rate of each nucleotide site is computed by looking at the frequency with which sequence pairs differ at that site as a function of their evolutionary distance. Variability maps constructed by this method provide a much more accurate and objective description of site-to-site variability than visual inspection of sequence alignments.

INTRODUCTION

A few years ago, quantitative substitution rate maps of the 5S rRNA, the small subunit ribosomal RNA (SSU rRNA), and the large subunit ribosomal RNA (LSU rRNA) of bacteria (1) and of the SSU rRNA of eukaryotes (2) were published. These maps were constructed by applying the substitution rate calibration (SRC) method that defines the variability or substitution rate of each nucleotide site as its evolutionary rate relative to the average evolutionary rate of all the nucleotide sites of the molecule (3,4). The variability maps constructed in this way clearly showed the distribution of variable and conserved positions in the different rRNA molecules.

Until recently, a reliable variability map for the eukaryotic LSU rRNA could not be constructed because too few complete sequences were available. However, due to recent sequencing efforts, the number of complete eukaryotic LSU rRNA sequences has increased significantly (5,6, unpublished results) and currently about 80 complete sequences are available for organisms belonging to the so-called crown of evolution (7), which now allows the construction of a reliable and detailed map showing the substitution rates in the LSU rRNA of eukaryotes.

Detailed information about the variability or conservation of nucleotide positions in ribosomal RNA is important for several reasons. The variability maps can be interpreted in terms of higher order structure and function and they facilitate the selection of areas suitable for the design of PCR primers and

hybridization probes. In addition, the measurement of site variability is important from a phylogenetic point of view. While conserved areas can be used to unravel old relationships, the more variable regions can be used to study evolution between closely related organisms (e.g. 8). Regarding phylogenetic tree construction, the study of site variability in molecules has gained much interest lately. Newly developed tree construction methods take into account differences in nucleotide substitution rates, which leads to more consistent tree topologies (9–11).

MATERIALS AND METHODS

Sequence alignment and estimation of substitution rates

The LSU rRNA database, established at the University of Antwerp (UIA) in 1994, is continuously updated by scanning the international nucleotide sequence libraries such as GenBank and EMBL for corrected or new ribosomal RNA sequences. In general, only complete or nearly complete sequences are compiled. All ribosomal RNA sequences are stored in the form of an alignment that is based on the secondary structure adopted for the molecule (e.g. 12–14). All sequences in the database are aligned with the DCSE sequence editor (15). Beside the primary and secondary structure information, literature references, accession numbers and detailed taxonomic information about the organism from which the sequence was derived are also compiled. For more information about the LSU rRNA database and its contents, we refer to the latest database issue of Nucleic Acids Research (16). The easiest way to obtain data is through the World Wide Web at URL <http://rrna.uia.ac.be/lsu/>

Estimation of substitution rates and construction of a variability map was done as described previously (1,3,4). In short, substitution rates are estimated by looking at the frequency with which sequence pairs differ at each site as a function of the distance between the sequence pairs (3). Substitution rates or variabilities v are estimated for every site in the sequence alignment that is not absolutely conserved and contains a nucleotide in at least 25% of the aligned sequences. Then, after estimation of all substitution rates, alignment positions are grouped into sets of similar rate. A spectrum of relative nucleotide substitution rates is thus obtained (1,4). A color map, superimposed on the secondary structure model of the LSU rRNA, can then be constructed by dividing the nucleotides into

*To whom correspondence should be addressed at present address: Department of Biology, University of Konstanz, D-78457 Konstanz, Germany.
Tel: +49 7531 88 2763; Fax: +49 7531 88 3018; Email: yves.vandeppeer@uni-konstanz.de

Table 1. Eukaryotic LSU rRNA sequences (and accession numbers) used for the nucleotide substitution rate calibration

Animals (18)		<i>Lithophragma trifoliata</i>	AF036501
<i>Acipenser brevirostrum</i>	U34340	<i>Mitella pentandra</i>	AF036502
<i>Aedes albopictus</i>	L22060	<i>Oryza sativa</i>	M11585, M16845
<i>Anguilla rostrata</i>	U34342	<i>Parnassia fimbriata</i>	AF036496
<i>Anopheles albimanus</i>	L78065	<i>Peltoboykinia tellimoides</i>	AF036499
<i>Caenorhabditis elegans</i>	X03680	<i>Plumbago auriculata</i>	AF036492
<i>Drosophila melanogaster</i>	M29800	<i>Saxifraga mertensiana</i>	AF036498
<i>Dugesia tigrina</i>	U78718	<i>Sinapis alba</i>	X66325
<i>Herdmania momus</i>	X53538	<i>Tellima grandiflora</i>	AF036500
<i>Homo sapiens</i>	M11167, J01866	<i>Tragopogon dubius</i>	AF036493
<i>Latimeria chalumnae</i>	U34336	Green algae (1)	
<i>Lepidosiren paradoxa</i>	U34337	<i>Chlorella ellipsoidea</i>	D17810
<i>Mus musculus</i>	X00525, J00623	Heterokont algae and relatives (8)	
<i>Neoceratodus forsteri</i>	U34338	<i>Hypochytrium catenoides</i>	X80345, X80346
<i>Oncorhynchus mykiss</i>	U34341	<i>Nannochloropsis salina</i>	Y07975, Y07974
<i>Protopterus aethiopicus</i>	U34339	<i>Ochromonas danica</i>	Y07977, Y07976
<i>Rattus norvegicus</i>	V01270	<i>Phytophthora megasperma</i>	X75631, X75632
<i>Xenopus borealis</i>	X59733	<i>Prorocentrum micans</i>	X16108, M14649
<i>Xenopus laevis</i>	X02995	<i>Scytosiphon lomentaria</i>	D16558
Fungi (11)		<i>Skeletonema pseudocostatum</i>	Y11512, Y11511
<i>Arxula adeninivorans</i>	Z50840	<i>Tribonema aequale</i>	Y07979, Y07978
<i>Blastocladiella emersonii</i>	X90411, X90410	Apicomplexans (6)	
<i>Candida albicans</i>	X70659, L07796	<i>Cryptosporidium parvum</i>	AF040725
<i>Cryptococcus neoformans</i>	L14067	<i>Eimeria tenella</i>	AF026388
<i>Entomophaga aulicae</i>	U35394	<i>Neospora caninum</i>	AF001946
<i>Pneumocystis carinii</i>	M86760	<i>Plasmodium falciparum</i>	U21939
<i>Saccharomyces cerevisiae</i>	J01355, K01048	<i>Theileria parva</i>	AF013419
<i>Saccharomycopsis fibuligera</i>	U09238, U10409	<i>Toxoplasma gondii</i>	X75429
<i>Schizosaccharomyces japonicus</i>	Z32848	Ciliates (4)	
<i>Schizosaccharomyces pombe</i>	Z19578	<i>Tetrahymena pyriformis</i>	X54004, M10752
<i>Tricholoma matsutake</i>	U62964	<i>Tetrahymena thermophila</i>	X54512
Land plants (21)		<i>Spathidium amphoriforme</i>	Unpublished
<i>Acorus gramineus</i>	AF036490	<i>Euplotes aediculatus</i>	Unpublished
<i>Arabidopsis thaliana</i>	X52320	Haptophytes (2)	
<i>Brassica napus</i>	D10840	<i>Phaeocystis antarctica</i>	Unpublished
<i>Drimys winteri</i>	AF036491	<i>Prymnesium patelliferum</i>	Unpublished
<i>Ephedra distachya</i>	AF036489	Red algae (2)	
<i>Eucryphia lucida</i>	AF036494	<i>Gracilaria verrucosa</i>	Y11508, Y11507
<i>Fragaria ananassa</i>	X58118, X15589	<i>Palmaria palmata</i>	Y11506
<i>Funaria hygrometrica</i>	X99331, X74114	Other (4)	
<i>Gnetum gnemon</i>	AF036488	<i>Chlorarachnion sp.</i>	Unpublished
<i>Hamamelis virginiana</i>	AF036495	<i>Pedinomonas minutissima</i>	U58510
<i>Jepsonia parryi</i>	AF036497	<i>Guillardia theta</i>	Unpublished
		<i>Guillardia theta nucleomorph</i>	Y11510, Y11509

Duplicate sequences belonging to the same species were omitted from the analysis. When two accession numbers are given, the first one is that of the 28S rRNA, the second one that of the 5.8S rRNA.

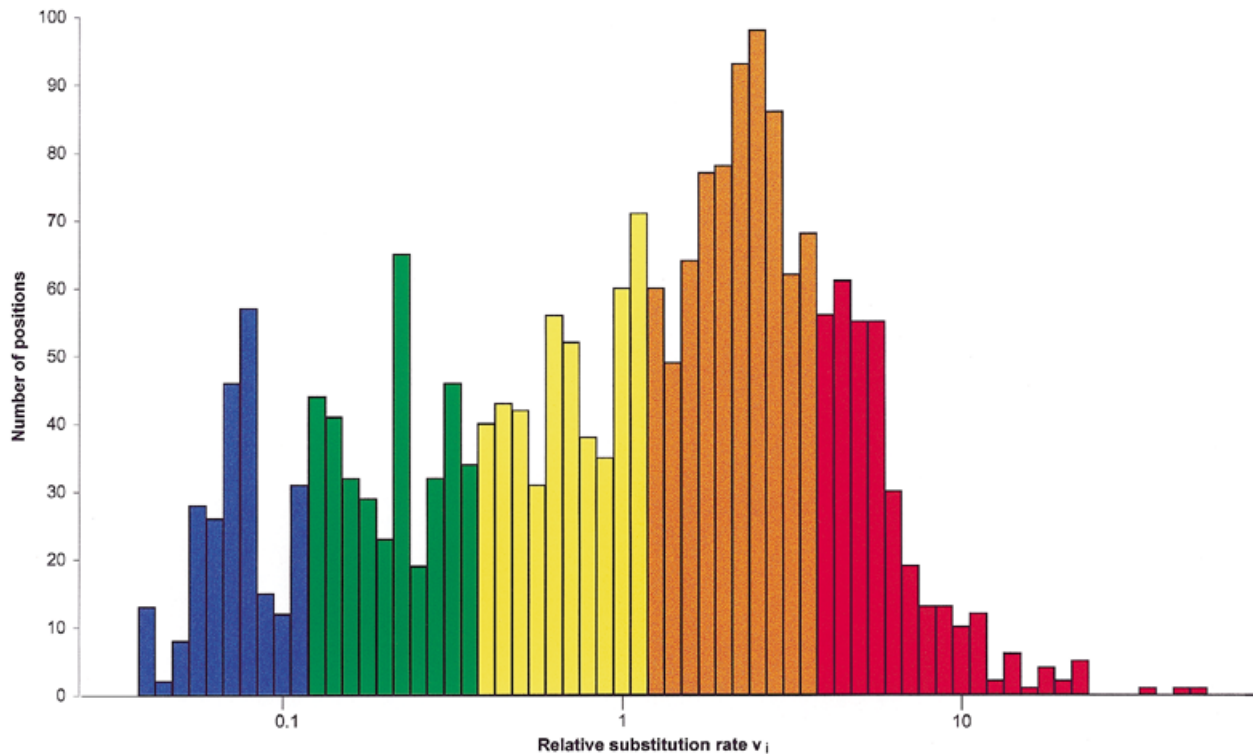


Figure 1. Distribution of the relative substitution rates, estimated from an alignment of 77 eukaryotic LSU rRNAs. The species used are specified in Table 1. Rates were estimated for 2153 alignment positions. Not included are 923 positions that are identical in all known sequences and 474 positions that contain a nucleotide in <25% of the aligned sequences. Sets of nucleotides are indicated in the same colors as used on the variability map of Figure 3.

different variability subsets and assigning a different color to each of these.

Once the shape of the spectrum is known, it is also possible to derive a new equation describing the evolutionary distance between two sequences as a function of the observed number of differences, i.e. the dissimilarity. This new equation can then be applied to tree construction, and several successful applications taking into account among-site rate variation in ribosomal RNAs have been described elsewhere (e.g. 10,17).

For the present study, 77 complete LSU rRNA sequences were analyzed, listed in Table 1. Duplicate sequences belonging to the same species were omitted from the analysis. The color map for eukaryotic LSU rRNA presented in this study can also be consulted on-line at URL <http://bioc-www.uia.ac.be/u/yvdp/>. Color maps of bacterial 5S rRNA, bacterial SSU rRNA, bacterial LSU rRNA and eukaryotic SSU rRNA, published previously (1,2), can be found there as well. Secondary structures were drawn with the software tool RNAviz (18).

RESULTS AND DISCUSSION

Variability map of eukaryotic LSU rRNA

The substitution rate spectrum obtained for eukaryotic LSU rRNA is shown in Figure 1. A different color is assigned to each of the different subsets. The relative rate limits of the subsets, and the corresponding colors used in the variability map, are as follows:

$0 < v_i < 10^{-0.925}$ (blue); $10^{-0.925} < v_i < 10^{-0.425}$ (green); $10^{-0.425} < v_i < 10^{+0.075}$ (yellow); $10^{+0.075} < v_i < 10^{+0.575}$ (orange); $v_i \geq 10^{+0.575}$ (red). Since the rate distribution is not rectangular (Fig. 1), some colors are more abundant than others.

Figure 2 shows the secondary structure model of the LSU rRNA of the yeast *Saccharomyces cerevisiae* while Figure 3 shows the variability of the nucleotide sites of LSU rRNA mapped in the same shape. Colors attributed to different sites are as described above. Absolutely conserved (invariant) positions ($v_i = 0$) are indicated in purple while sites colored in pink belong to areas that are very variable, but that are deleted in too many sequences to allow a sufficiently accurate measurement of their relative evolutionary rate. The color map for LSU rRNA gives a much more detailed and quantitative description of positional variability than the crude distinction between variable and conserved areas that is often made by visual inspection of sequence alignments.

It can be seen that in general the two nucleotides of a base pair have the same or a neighboring color, i.e. they are about equally variable. This is as expected, since the substitution of a base-paired nucleotide generally requires a compensating substitution in the opposite strand. Exceptions are mostly due to the fact that in some cases a particular base, usually a G or a U, seems to be required in one strand, but the existence of G-U pairs aside from G-C and A-U pairs allows the complementary base to change more freely.

As can be seen in Figure 3, 10 highly variable areas (orange, red and pink) can be distinguished in the eukaryotic LSU

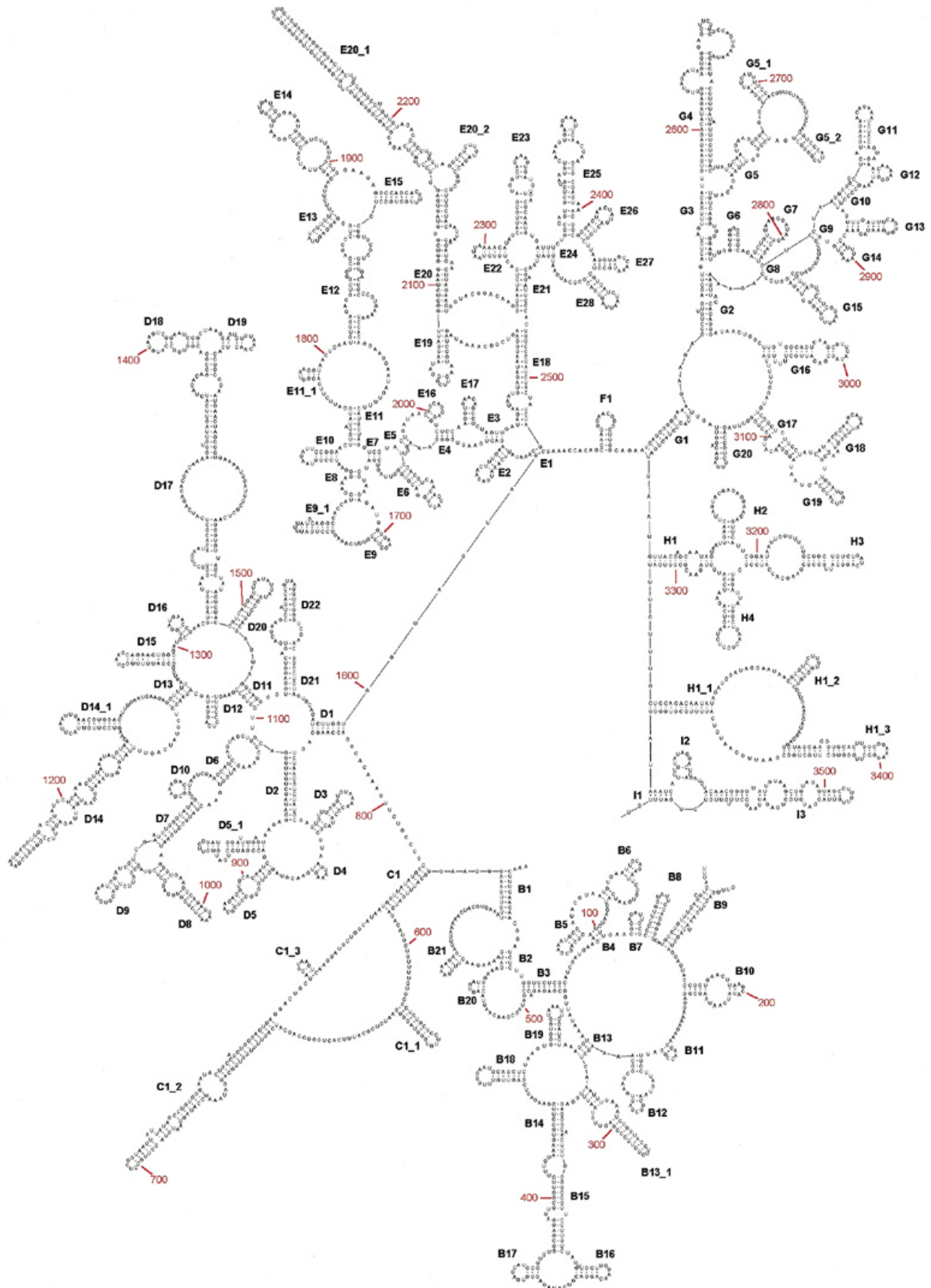


Figure 2. Secondary structure model for the LSU rRNA of *S.cerevisiae*. The sequence is written clockwise from 5' to 3', sites are numbered in red every 100 nucleotides. Helix numbering is according to De Rijk *et al.* (16).

Table 2. Fraction of more conserved and less conserved sites in double-stranded areas and in different types of loops

Relative substitution rate v_i	No. of sites in different structural elements					
	multi-branched loops	hairpin loops	internal- and bulge loops ^a	all single stranded elements	helices	entire molecule ^b
$v_i < 10^{-0.925}$	252	180	181	613	548	1161
$v_i \geq 10^{-0.925}$	322	201	268	791	1124	1915
Total	574	381	449	1404	1672	3076
% with $v_i < 10^{-0.925}$	43.9	47.2	40.3	43.7	32.8	37.7
% with $v_i \geq 10^{-0.925}$	56.1	52.8	59.7	56.3	67.2	62.3

^aThe number of sites in bulge loops was too small (43 sites) to make a separate calculation meaningful.

^bThis includes all sites for which a relative substitution rate was computed, i.e. all but those deleted in >75% of the sequences and indicated in pink on the color map of Figure 3.

rRNA. These are formed by the following helices: B8; B13_1 to B16; the whole area denoted as C; D5 and D5_1; D14_1; D20; E20_1 to E20_2; G5_1 to G5_2; and H1_1 to H1_3. Many of the variable areas are characterized by major size variations. For example, the areas enclosed by helices C1 and E20, the entire area H1_n, and to a lesser extent helix D20, are also hot-spots for extremely variable insertions. These insertions were first described by Hassouna *et al.* (19) who referred to them as D(ivergent)-domains. As a rule, strong length heterogeneity seems to be most common in apical helices, i.e. those ending in a hairpin loop. Helices formed by long distance interactions, i.e. those bounded by multibranching loops, have less freedom to change in length. It should also be noted that the LSU rRNA molecule contains a number of potential branching points that bear additional helices in a limited set of species. For example, helices B14 and B15, though separated only by an internal loop in *S.cerevisiae*, were numbered differently because a potential branching point separates them.

Beside highly variable regions, several regions of conserved nature can be distinguished in the LSU rRNA. For the LSU rRNA of prokaryotes and for the SSU rRNA of prokaryotes and eukaryotes a stronger sequence conservation in single-strands than in double-strands has been reported (1,8,20–22). In order to examine this quantitatively in the case of eukaryotic LSU rRNA, separate substitution rate spectra were measured for nucleotide sites involved in base pairing and for those forming part of each type of single-stranded structural elements: multibranching-, hairpin-, internal- and bulge-loops. These spectra (not shown) were not very different from that measured for the entire molecule (Fig. 1). However, it should be remembered that the spectra show only the distribution of sites with a measurable substitution rate and do not include sites that are identical in all hitherto sequenced molecules. The latter have $v_i = 0$ and therefore cannot be represented on a logarithmic scale. In order to demonstrate the greater tendency of conservation in single-

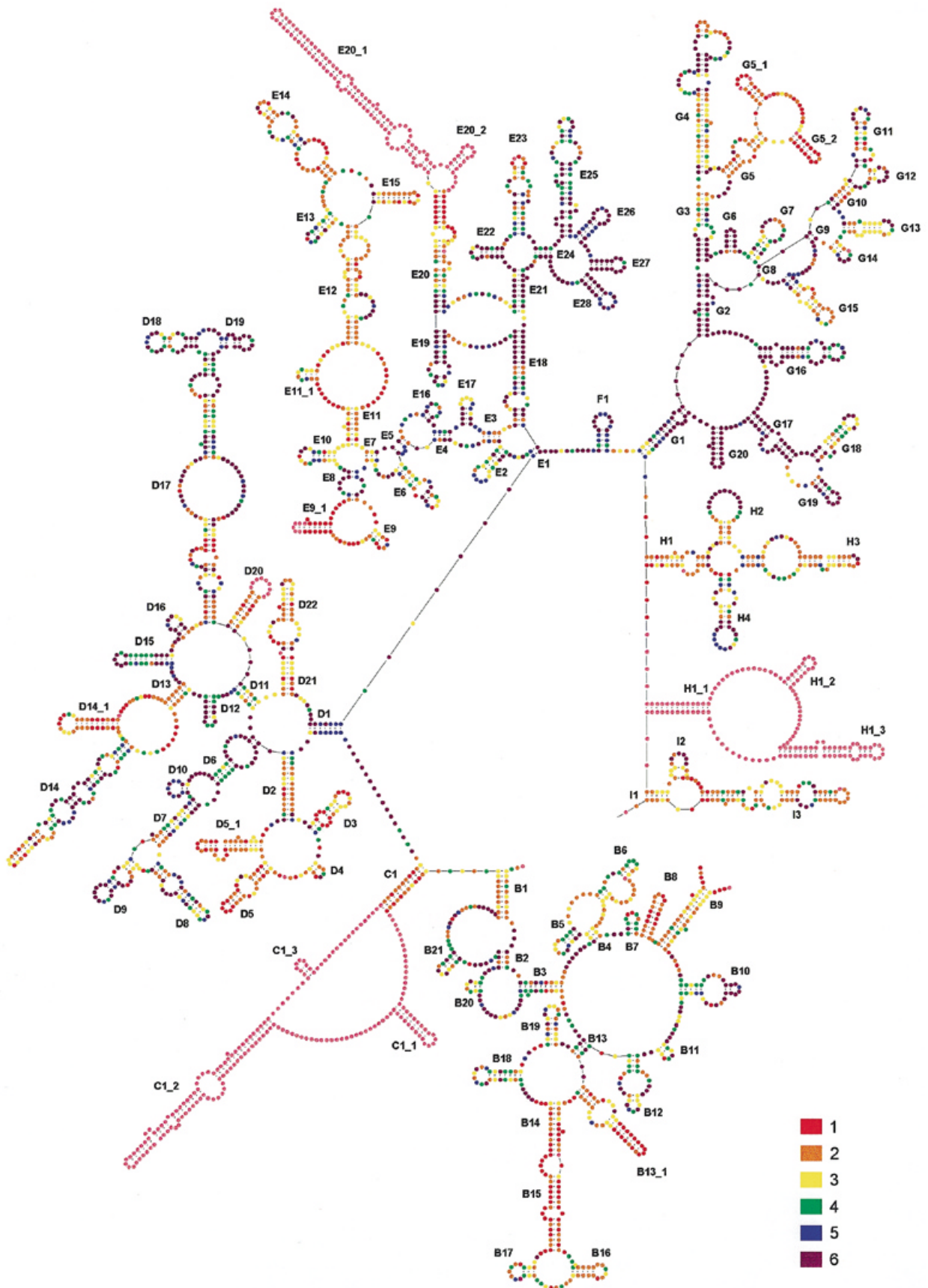
stranded areas, sites were divided arbitrarily in more conserved ones ($v_i < 10^{-0.925}$, i.e. purple and blue in Figs 1 and 3) and more variable ones ($v_i \geq 10^{-0.925}$, i.e. green to red in Figs 1 and 3). Table 2 shows that the fraction of more conserved sites is larger in single-stranded structural elements than in helices. In Figure 4 the fraction of more conserved sites ($v_i < 10^{-0.925}$) is calculated in both single- and double-stranded regions over the entire molecule in a sliding window of 50 nucleotides. In this calculation positions that contain a nucleotide in <25% of the sequences (pink in Fig. 3) were included with the most variable ones ($v_i \geq 10^{-0.925}$). This graph shows more clearly that, overall, single-stranded regions have a higher fraction of conserved positions than double-stranded regions of the LSU rRNA molecule. One notable exception is the area between positions 2000 and 2400, which corresponds approximately to the region covered by helices E16–E25.

Very often, specific functions can be ascribed to regions of the molecule that are conserved in structure (e.g. 23,24). One of the most conserved areas, both in bacterial (1) and eukaryotic LSU rRNA is the large multibranching loop in area G and the helices surrounding it (Fig. 3). This structure is generally considered to be the major element of the peptidyl transferase center of the ribosome, which is the catalytic center responsible for the peptide bond formation (25,26). Other highly conserved structures in the LSU rRNA are helices D18–D19 which are part of the so-called GTPase center of the ribosome, helices E21–E28 and the hairpin loop of helix H2 (25).

ACKNOWLEDGEMENTS

We want to thank Linda Medlin for DNA of haptophytes. Our research is supported by the Special Research Fund of the University of Antwerp and by the Fund for Scientific Research, Flanders. Yves Van de Peer is Research Fellow of the Fund for Scientific Research, Flanders.

Figure 3. (Next page) Variability map superimposed on the LSU rRNA secondary structure model of *S.cerevisiae*. Nucleotides are subdivided into five groups of increasing variability (see text for details). The most variable positions are in red, the least variable ones in blue. Absolutely conserved positions in all structures hitherto known are indicated in purple. Hypervariable regions that were not taken into consideration for rate calibration, because they are absent in >75% of the eukaryotic sequences considered, are indicated in pink. These include C1_1 to C1_3; E20_1 to E20_2; H1_1 to H1_3, the hairpin loops of D20 and E9_1, as well as individual nucleotides peculiar to the *S.cerevisiae* LSU rRNA.



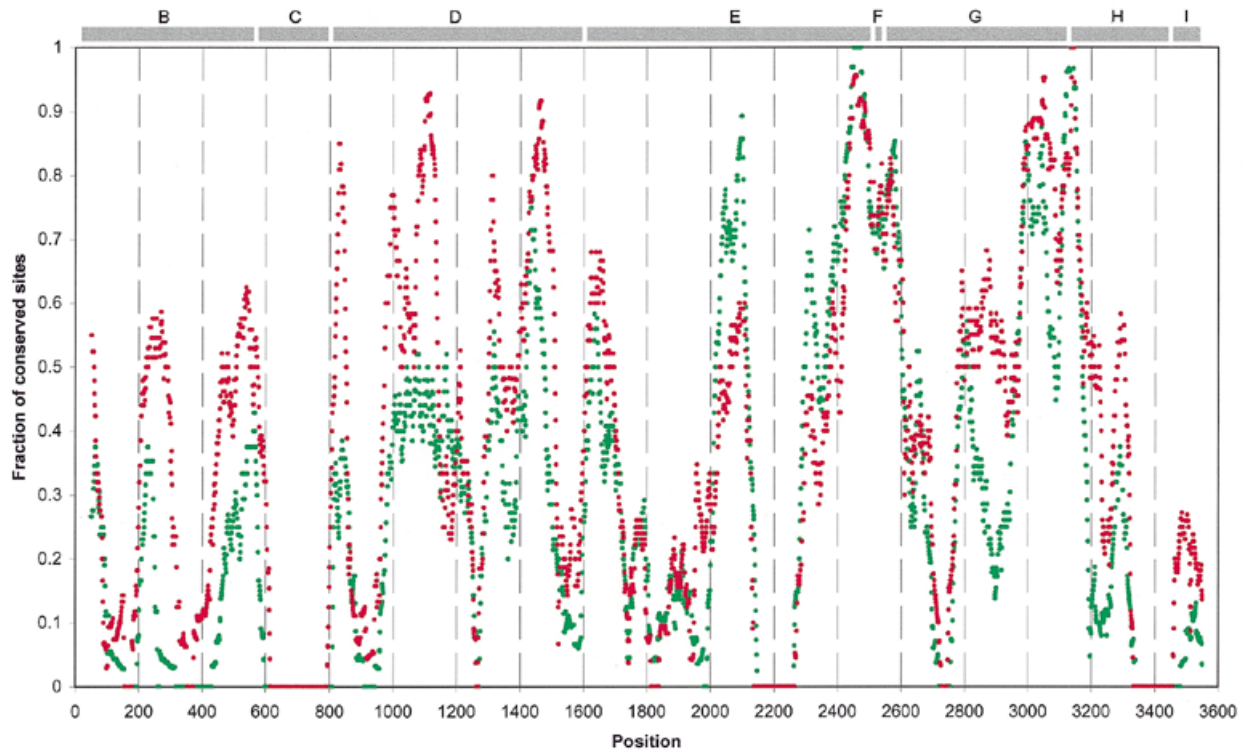


Figure 4. Graph showing the fraction of conserved sites, counted with a sliding window of 50 nucleotides. Each red dot represents the fraction of positions with $v_i < 10^{-0.925}$ in single stranded sites, each green dot represents this fraction in helical sites. The location of structural areas B to I is indicated by bars on top of the figure. Long areas with a fraction of conserved sites equal to 0 around positions 700, 2200 and 3400 consist of sites of presumably high but unmeasurable variability (pink dots in Fig. 3).

REFERENCES

1. Van de Peer, Y., Chapelle, S. and De Wachter, R. (1996) *Nucleic Acids Res.*, **24**, 3381–3391.
2. Van de Peer, Y., Jansen, J., De Rijk, P. and De Wachter, R. (1997) *Nucleic Acids Res.*, **25**, 111–116.
3. Van de Peer, Y., Neefs, J.-M., De Rijk, P. and De Wachter, R. (1993) *J. Mol. Evol.*, **37**, 221–232.
4. Van de Peer, Y., Van der Auwera, G. and De Wachter, R. (1996) *J. Mol. Evol.*, **42**, 201–210.
5. Van der Auwera, G. and De Wachter, R. (1997) *J. Mol. Evol.*, **45**, 84–90.
6. Van der Auwera, G., Hofmann, C.J., De Rijk, P. and De Wachter, R. (1998) *Mol. Phyl. Evol.*, **10**, 333–342.
7. Knoll, A.H. (1992) *Science*, **256**, 622–627.
8. Abouheif, E., Zardoya, R. and Meyer, A. (1998) *J. Mol. Evol.*, **47**, 394–405.
9. Olsen, G.J. (1987) *Cold Spring Harbor Symp. Quant. Biol.*, **LII**, 825–837.
10. Van de Peer, Y., Rensing, S., Maier, U.-G. and De Wachter, R. (1996) *Proc. Natl Acad. Sci. USA*, **93**, 7732–7736.
11. Yang, Z. (1996) *Trends Ecol. Evol.*, **11**, 367–372.
12. Gutell, R.R., Gray, M.W. and Schnare, M.N. (1993) *Nucleic Acids Res.*, **21**, 3055–3074.
13. De Rijk, P., Van de Peer, Y., Chapelle, S. and De Wachter, R. (1994) *Nucleic Acids Res.*, **22**, 3495–3501.
14. De Rijk, P., Van de Peer, Y. and De Wachter, R. (1996) *Nucleic Acids Res.*, **24**, 92–97.
15. De Rijk, P. and De Wachter, R. (1993) *Comput. Appl. Biosci.*, **9**, 735–740.
16. De Rijk, P., Robbrecht, E., de Hoog, S., Caers, A., Van de Peer, Y., De Wachter, R. (1999) *Nucleic Acids Res.*, **27**, 174–178.
17. Van de Peer, Y. and De Wachter, R. (1997) *J. Mol. Evol.*, **45**, 619–630.
18. De Rijk, P. and De Wachter, R. (1997) *Nucleic Acids Res.*, **25**, 4679–4684.
19. Hassouna, N., Michot, B. and Bachelier, J.-P. (1984) *Nucleic Acids Res.*, **12**, 3563–3581.
20. Vawter, L. and Brown, W.M. (1993) *Genetics*, **134**, 597–608.
21. Rzhetsky, A. (1995) *Genetics*, **141**, 771–783.
22. Otsuka, J., Terai, G. and Nakano, T. (1999) *J. Mol. Evol.*, **48**, 218–235.
23. Egebjerg, J., Larsen, N. and Garrett, R.A. (1990) In Hill, W.E., Dahlberg, A., Garrett, R.A., Moore, P.B., Schlessinger, D. and Warner, J.R. (eds), *The Ribosome. Structure, Function and Evolution*. American Society of Microbiology, Washington, DC, pp. 168–179.
24. Noller, H.F., Moazed, D., Stern, S., Powers, T., Allen, P.N., Robertson, J.M., Weiser, B. and Triman, K. (1990) In Hill, W.E., Dahlberg, A., Garrett, R.A., Moore, P.B., Schlessinger, D. and Warner, J.R. (eds), *The Ribosome. Structure, Function and Evolution*. American Society of Microbiology, Washington, DC, pp. 73–92.
25. Raué, H.A., Musters, W., Rutgers, C.A., van't Riet, J. and Planta, R.J. (1990) In Hill, W.E., Dahlberg, A., Garrett, R.A., Moore, P.B., Schlessinger, D. and Warner, J.R. (eds), *The Ribosome. Structure, Function and Evolution*. American Society of Microbiology, Washington, DC, pp. 217–235.
26. Bocchetta, M., Xiong, L. and Mankin, A.S. (1998) *Proc. Natl Acad. Sci. USA*, **95**, 3525–3530.